

Enhancing animal breeding through quality control in genomic data - a review

Jungjae Lee^{1#}, Jong Hyun Jung^{2#}, Sang-Hyon Oh^{3*}

¹Jenomics Jenetics Company, Pyeongtaek 17869, Korea

²Jung P&C Institute, Yongin 16950, Korea

³Division of Animal Science, Gyeongsang National University, Jinju 52725, Korea



Received: Sep 14, 2024

Revised: Sep 30, 2024

Accepted: Oct 1, 2024

#These authors contributed equally to this work.

*Corresponding author

Sang-Hyon Oh

Division of Animal Science,
Gyeongsang National University, Jinju
52725, Korea.

Tel: +82-55-772-3285

E-mail: shoh@gnu.ac.kr

Copyright © 2024 Korean Society of
Animal Sciences and Technology.

This is an Open Access article
distributed under the terms of the
Creative Commons Attribution
Non-Commercial License (<http://creativecommons.org/licenses/by-nc/4.0/>) which permits unrestricted
non-commercial use, distribution, and
reproduction in any medium, provided
the original work is properly cited.

ORCID

Jungjae Lee

<https://orcid.org/0000-0002-6145-8862>

Jong Hyun Jung

<https://orcid.org/0000-0003-3667-7710>

Sang-Hyon Oh

<https://orcid.org/0000-0002-9696-9638>

Competing interests

No potential conflict of interest relevant
to this article was reported.

Funding sources

Not applicable.

Acknowledgements

This work was carried out with the
support of "Cooperative Research
Program for Agriculture Science and
Technology Development (Project No.

Abstract

High-throughput genotyping and sequencing has revolutionized animal breeding by providing access to vast amounts of genomic data to facilitate precise selection for desirable traits. This shift from traditional methods to genomic selection provides dense marker information for predicting genetic variants. However, the success of genomic selection heavily depends on the accuracy and quality of the genomic data. Inaccurate or low-quality data can lead to flawed predictions, compromising breeding programs and reducing genetic gains. Therefore, stringent quality control (QC) measures are essential at every stage of data processing. QC in genomic data involves managing single nucleotide polymorphism (SNP) quality, assessing call rates, and filtering based on minor allele frequency (MAF) and Hardy-Weinberg equilibrium (HWE). High-quality SNP data is crucial because genotyping errors can bias the estimates of breeding values. Cost-effective low-density genotyping platforms often require imputation to deduce missing genotypes. QC is vital for genomic selection, genome-wide association studies (GWAS), and population genetics analyses because it ensures data accuracy and reliability. This paper reviews QC strategies for genomic data and emphasizes their applications in animal breeding programs. By examining various QC tools and methods, this review highlights the importance of data integrity in achieving successful outcomes in genomic selection, GWAS, and population analyses. Furthermore, this review covers the critical role of robust QC measures in enhancing the reliability of genomic predictions and advancing animal breeding practices.

Keywords: Animal breeding, Genomic selection, Quality control, Single nucleotide polymorphism, Genome-wide association studies

INTRODUCTION

The rapid evolution of genomic technologies has transformed the landscape of animal breeding. High-throughput genotyping and sequencing provides breeders with access to vast amounts of genomic data and enables the precise selection of desirable traits [1]. These advancements have shifted traditional breeding methods to genomic selection, which leverages dense marker information to predict the genetic variants of individuals [2]. However, the success of genomic selection depends heavily on the

RS-2023-00232087)* Rural Development Administration, Korea.

Availability of data and material

Upon reasonable request, the datasets of this study can be available from the corresponding author.

Authors' contributions

Conceptualization: Lee JJ, Jung JH, Oh SH.
Data curation: Lee JJ, Jung JH, Oh SH.
Formal analysis: Lee JJ, Jung JH, Oh SH.
Methodology: Lee JJ, Jung JH, Oh SH.
Validation: Lee JJ, Jung JH, Oh SH.
Investigation: Lee JJ, Jung JH, Oh SH.
Writing - original draft: Lee JJ, Jung JH, Oh SH.
Writing - review & editing: Lee JJ, Jung JH, Oh SH.

Ethics approval and consent to participate

This article does not require IRB/ACUC approval because there are no human and animal participants.

accuracy and quality of the genomic data. Inaccurate or low-quality data can lead to inaccurate predictions that can compromise breeding programs and reduce their genetic gains [3]. Therefore, to ensure reliable predictions and maximize the potential of genomic selection, it is essential to implement stringent quality control (QC) measures at every stage of data processing.

Genomic data QC has several key components including the management of single nucleotide polymorphism (SNP) quality, the assessment of call rates, and filtering based on minor allele frequency (MAF) and Hardy-Weinberg equilibrium (HWE) [4]. High-quality SNP data is indispensable because errors in genotyping can lead to biased estimates of breeding values, which decreases the effectiveness of selection strategies [5]. Moreover, cost-effective low-density genotyping platforms often suffer from incomplete marker data so it is necessary to use imputation to deduce the missing genotypes [6].

QC processes are crucial for genomic selection, genome-wide association studies (GWAS), and population genetics analyses. These processes help ensure that the genomic data is accurate, reliable, and free from biases introduced by genotyping errors, population stratification, or other confounding factors [7,8]. This paper reviews QC strategies for genomic data and their applications in animal breeding programs. By examining various QC tools and methods, this paper aims to show the critical role that data integrity plays in achieving successful outcomes in genomic selection, GWAS, and population analyses [4,5].

GENOTYPING METHODS

Whole-genome sequencing (WGS)

WGS is a comprehensive method for analyzing the entire genome. Due to the decreased cost of sequencing and the ability to produce large amounts of genomic data, WGS has become a powerful tool for genomic research. SNP calling from WGS genomic data involves a series of critical steps to ensure accurate identification of genetic variants. The process starts with raw data preprocessing, where tools like FastQC evaluate the read quality [9]. This step is followed by trimming to remove adapters and low-quality bases by using either Trimmomatic or Cutadapt [10,11].

The cleaned reads are then aligned to a reference genome with BWA-MEM or Bowtie2 to generate SAM/BAM files [12,13]. These files are subsequently sorted, indexed, and processed to mark polymerase chain reaction (PCR) duplicates with Samtools, while the base quality scores are recalibrated using GATK [14,15]. Variant calling is performed using tools such as GATK's HaplotypeCaller, FreeBayes, or Bcftools, which identify SNPs based on differences between the sequenced reads and the reference genome [15–17].

In post-calling, variants undergo filtering to remove false positives via GATK's hard filtering or Variant Quality Score Recalibration (VQSQR). The filtered SNPs are then annotated with functional information using tools like ANNOVAR or SnpEff [18,19]. Quality checks include the use of VCFtools for statistical analysis and IGV for visualization, and ensure the reliability of the called SNPs [16,20]. Joint genotyping across multiple samples and using population-specific reference panels are recommended to enhance the accuracy of SNP calling in WGS.

SNP arrays

SNP arrays have significantly advanced genomic research in animal science by enabling the large-scale genotyping of SNPs. The development of SNP arrays began in the early 2000s to meet the demand for efficient and cost-effective methods to genotype large numbers of SNPs across the genome [21,22]. Early arrays marked a significant advancement by allowing simultaneous genotyping of thousands of SNPs, facilitating GWAS and the study of genetic variation in

populations [21].

Over time, these arrays have evolved to include higher-density SNPs to improve coverage and accuracy, as seen in the Illumina BovineSNP50 array which has become a standard tool in cattle genomics [23,24]. Today, SNP arrays are essential for selecting desirable traits, estimating genetic merit, and managing inbreeding in animal breeding [1,2]. QC of SNP array data is crucial for ensuring accurate and reliable results, and involves assessing call rates, filtering based on MAF, and checking for HWE [4]. Tools such as PLINK and GenomeStudio are commonly used in these QC processes [5,25].

QC IN ANIMAL GENOMICS

Minor allele frequency (MAF)

MAF is a key metric in genetic studies. It represents the frequency at which the less common allele occurs in a given population. MAF is important for identifying rare variants which may not significantly contribute to overall genetic variation but can be crucial in specific contexts. MAF is calculated by determining the frequency of both alleles at a locus and taking the minimum of these two values. For example, if allele A has a frequency of 0.8 and allele a has a frequency of 0.2, the MAF would be 0.2. SNPs with very low MAFs, typically below 0.01 or 0.05, are often excluded from analyses because they may represent sequencing errors or lack statistical power in association studies [5].

Tools like PLINK and VCFtools [5,16] are widely used to calculate MAF, with PLINK's `--freq` command being particularly popular [4]. In animal breeding, many researchers set threshold values for MAF to balance the need for sufficient variation while minimizing noise from rare variants. Typically, MAF thresholds in animal breeding studies range from 0.01 to 0.05 depending on the study's objectives and the population structure being analyzed. For instance, a study on dairy cattle by Pryce et al. [26] and Kim et al. [27] used a MAF threshold of 0.01 to ensure that the SNPs included were sufficiently informative for genomic predictions while also minimizing the influence of rare variants that might lead to spurious associations.

Call rate

Call rate is another critical QC metric that measures the proportion of successfully genotyped samples for a specific SNP. A high call rate indicates that a SNP has been consistently detected across the sample population, while a low call rate may suggest issues with the genotyping process, such as poor quality or technical errors [7].

The call rate is calculated by dividing the number of successful genotype calls for a SNP by the total number of samples, then multiplying by 100 to express it as a percentage.

$$\text{Call Rate} = \frac{\text{Number of successfully genotyped markers (or samples)}}{\text{Total number of markers (or samples)}} \times 100$$

For instance, if 95 out of 100 samples have a successful genotype call for a SNP, the call rate would be 95% [4]. Normally, markers with a call rate less than 95% are removed, though other studies have set more stringent or lenient thresholds depending on the study design and objectives. For example, some studies have removed markers with a call rate below 99% to ensure extremely high data quality [28], while others have used a more relaxed threshold of 90% when working with larger datasets [29].

Tools like PLINK, SNP & Variation Suite (SVS), and GenomeStudio are widely used for

calculating and filtering SNPs based on call rates because they offer robust functionalities for QC in genomic studies. PLINK is particularly popular due to its comprehensive command-line interface, where the `--missing` command calculates call rates at both the marker and sample levels, allowing researchers to easily filter out SNPs and samples that fall below the desired threshold [5]. SVS offers a user-friendly graphical interface and integrates various statistical tools, making it ideal for complex datasets and large-scale studies [30]. GenomeStudio by Illumina is another powerful tool specifically designed for managing and analyzing genotyping data with features for calculating call rates, identifying low-quality markers, and visualizing data for further inspection [25]. These tools are essential for ensuring that only high-quality data is used in subsequent analyses to improve the reliability of genomic outcomes.

Hardy-Weinberg equilibrium

HWE is a fundamental principle in population genetics. It states that allele and genotype frequencies in a population will remain constant from generation to generation in the absence of evolutionary influences [31]. Testing for HWE is an important QC step because deviations from this equilibrium can indicate issues such as genotyping errors, population stratification, or selection pressures [32]. To test for HWE, the observed genotype frequencies are compared to the expected frequencies under equilibrium conditions. For a biallelic SNP with alleles A and a, the expected genotype frequencies are p^2 for AA, $2pq$ for Aa, and q^2 for aa, where p and q represent the allele frequencies [33]. A chi-square test is commonly used to assess whether the differences between the observed and expected frequencies are statistically significant. Tools like PLINK and VCFtools are used to perform HWE tests [34]. SNPs that show significant deviation from HWE, typically with a p-value less than 0.001, are often excluded from analyses to prevent biases that could arise from genotyping errors or other confounding factors [4]. These QC metrics are foundational for ensuring high-quality genotypic data, forming the basis for accurate and reliable analyses in applications such as population analysis, GWAS, and genomic selection. Table 1 provides a summary of tools commonly used for QC steps, offering researchers practical options to streamline their workflows and enhance data integrity.

APPLICATION

Population analysis

Population analysis is invaluable for genomic studies in animal science because it enables researchers to assess the genetic structure, diversity, and evolutionary dynamics within and between populations. Accurately characterizing population structures is crucial for identifying subpopulations, measuring inbreeding levels, and understanding the genetic background of breeding populations, all of which are essential for maintaining genetic diversity and improving selection outcomes [35]. Tools such as PLINK, ADMIXTURE, and STRUCTURE are commonly employed to detect key characteristics for understanding the genetic landscape of animal populations, such as population stratification, admixture, and genetic differentiation [5,36]. For example, ADMIXTURE provides estimates of individual ancestry proportions. These estimates allow researchers to detect mixed genetic backgrounds that could influence trait analysis [36]. QC measures, such as filtering based on MAF, HWE, and genotyping call rates ensure the data used for population analysis is reliable [4,37]. MAF filtering helps exclude rare alleles that may introduce noise or result from genotyping errors [5]. Similarly, HWE filtering removes SNPs that deviate from expected frequencies due to selection or population substructures in order to prevent potential biases in the analysis [37]. Proper QC improves the accuracy of population structure analyses and mitigates the risk of

Table 1. Tool list for quality control processes

Tools	Function	Reference
GEMMA	Application of linear mixed models and related models to GWAS	[4]
PLINK	Run association analyses and perform QC and regression steps	[5]
FastQC	Quality control checks on raw sequence data	[9]
Trimmomatic	Trim and crop FASTQ data	[10]
Cutadapt	finds and removes adapter sequences, primers, poly-A tails	[11]
BWA-MEM	produce multiple primary alignments for different part of a query sequence	[12]
Bowtie2	aligning sequencing reads to long reference sequences	[13]
Samtools	Manipulate alignments in the SAM, BAM, and CRAM formats	[14]
GATK	Variant calling using sequencing data	[15]
VCFtools	Summarize, filter out, convert data into other file formats	[16]
FreeBayes	Bayesian genetic variant detector designed to find SNPs	[17]
SnEff	Annotation on genetic variants and predicts their effects on genes	[18]
ANNOVAR	Generate gene-based annotation	[19]
IGV	Visualization tool to simultaneously integrate and analyze multiple types of genomic data	[20]
GenomeStudio	Normalize, cluster, and call genotypes	[25]
SVS	Perform analyses and visualizations on genomic and phenotypic data	[33]
BEAGLE	Genotype calling, phasing, and genotype imputation	[39]
Fimpute	Haplotype estimation or phasing and genotype imputation	[40]
Impute2	Genotype imputation and haplotype phasing	[47]
Minimac	performs imputation with pre-phased haplotypes	[48]

confounding in subsequent analyses such as GWAS and genomic selection [4]. By accurately characterizing population structures, researchers can identify unique genetic markers and enhance their understanding of trait inheritance, and then design breeding strategies that optimize genetic gain and preserve diversity to support sustainable livestock production [35,36].

GWAS

GWAS are powerful tools for identifying genetic variants associated with complex traits in animal breeding such as growth traits, disease resistance, reproductive traits, and carcass traits [2,4]. The reliability of GWAS findings hinges on rigorous QC procedures that ensure high-quality data throughout the process. This begins with careful study design and population selection, where potential confounders like population stratification are addressed through methods such as Principal Component Analysis (PCA) and linear mixed models to correct for genetic structure within the population [38]. Phenotype data must be accurately collected and screened for outliers to minimize noise. Genotype data undergoes thorough QC, including filtering SNPs based on call rates, MAF, and deviations from HWE [4,5]. For instance, SNPs with low call rates are excluded to avoid unreliable data that could lead to false-positive associations, while MAF filtering focuses the analysis on common variants that are more likely to have sufficient statistical power to detect true associations. HWE filtering is employed to remove SNPs that significantly deviate from expected allele frequencies because such deviations may indicate genotyping errors or underlying selection pressures [5]. To reduce redundancy and computational burden, linkage disequilibrium (LD) pruning is performed and missing genotypes are often imputed via reference panels using Fimpute or BEAGLE [39,40]. Tools like PLINK and GEMMA are widely used to implement QC measures and conduct association tests because they offer a robust framework for analyzing large genomic datasets [4]. Statistical analysis in GWAS is carried out using models appropriate for

the trait under study, and corrections for multiple testing to mitigate the risk of false positives and meta-analysis may be employed when integrating results from multiple studies [41]. To ensure the robustness and high accuracy of the GWAS models, a 5-fold cross-validation is often used. In this method, the datasets are divided into five subsets. The model is iteratively trained on four subsets and tested on the remaining one to help validate the model's accuracy and mitigate overfitting [42]. The results from GWAS offer valuable genetic variants for traits which can be targeted in marker-assisted selection and genomic selection programs. Genomic selection aims to ultimately improve the genetic merit of livestock populations [2]. The Fig. 1 summarizes the genotype QC workflow, with an emphasis on data preparation, QC steps, and their applications.

Genomic selection

Genomic selection (GS) allows for the selection of animals based on SNP markers [43]. With the introduction of GS, animal breeding has dramatically advanced by overcoming the limitations of traditional selection methods like best linear unbiased prediction (BLUP) and marker-assisted selection [43,44]. GS relies on dense SNP data to estimate genomic breeding values, which are used to predict an individual's genetic potential for economically important traits [2]. The accuracy of GS models is dependent upon the quality of the genomic data and the reliability of GS models can be enhanced significantly by the inclusion of imputation methods to handle missing or low-density SNP data [45]. Imputation is beneficial in low-density platforms because it allows for the cost-effective use of genotyping while still leveraging the power of high-density SNP information. Imputation increases the accuracy of genomic predictions by inferring missing genotypes in order to improve the reliability of estimated breeding values even with fewer markers [6]. Several imputation tools, including FImpute [40], Beagle [39], Impute2 [46], and Minimac [47] are widely used in animal breeding to enhance the accuracy of GS models. Therefore, strict QC is essential [48]. QC methods, such as filtering SNPs based on call rates, MAF, and HWE, is critical to ensuring that the data is vigorous and reliable. High call rates are important because missing data can introduce

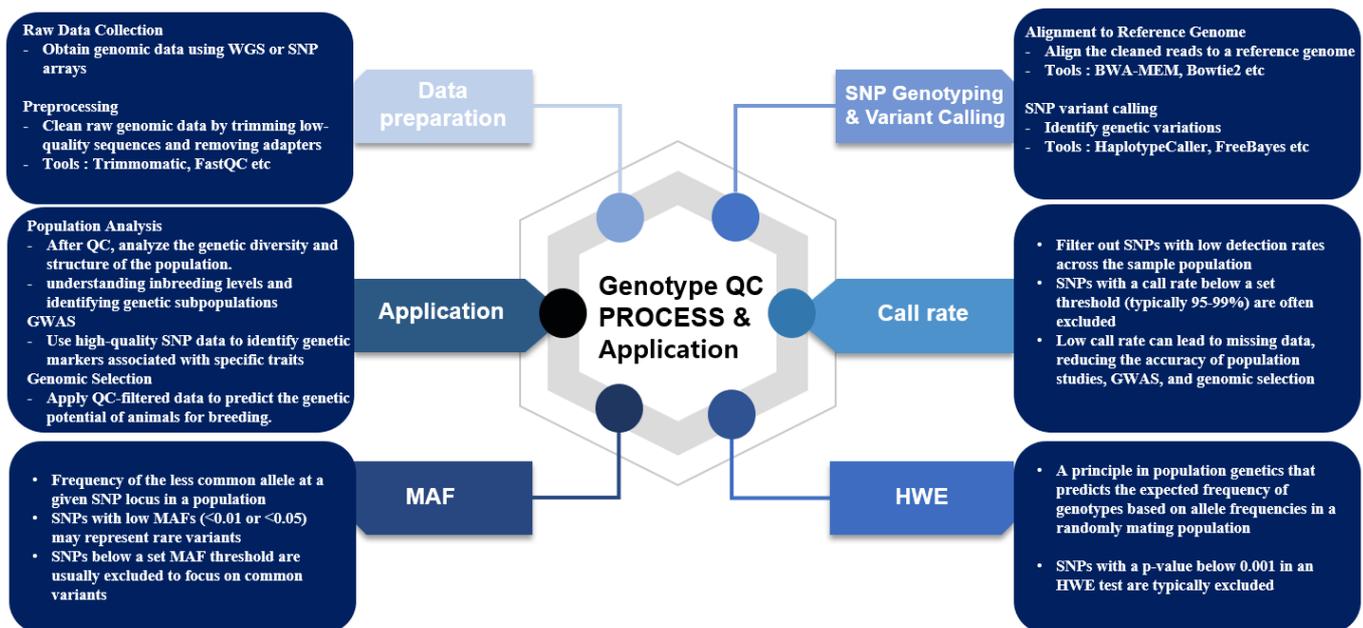


Fig. 1. Overall flowchart from data preparation to application in animal breeding.

bias and reduce the reliability of genomic estimated breeding values. Similarly, excluding SNPs with low MAF helps to avoid the noise associated with rare variants that may have little impact on prediction accuracy. Ensuring that SNPs conform to HWE expectations also prevents the inclusion of markers affected by selection, mutation, or other factors that could bias the GS models [4,5]. Advanced computational tools, such as genomic best linear unbiased prediction (GBLUP) and single-step BLUP (ssBLUP), and Bayesian methods (BayesA, BayesB, BayesC) integrate SNP effects across the genome to enhance the precision of breeding value predictions [49,50]. By using high-quality genomic data, GS enables breeders to make more accurate decisions that lead to faster genetic gains and the improvement of traits such as milk yield, growth rate, and carcass weight in livestock. This approach not only enhances the efficiency of breeding programs but also contributes to the long-term sustainability and productivity of animal populations [35].

CONCLUSION

High-throughput genotyping and sequencing has significantly advanced the field of animal breeding by enabling precise selection for desirable traits. However, the success of GS hinges on the accuracy and quality of the genomic data used. Rigorous QC measures are essential to ensure data integrity. These measures include SNP quality management, call rate assessment, and filtering based on MAF and HWE. These QC processes are crucial for GS, GWAS, and population genetics analyses. Implementing stringent QC strategies enhances the reliability of genomic predictions, which improves breeding programs and genetic gains. By maintaining high standards of data quality, researchers and breeders can make informed decisions that lead to sustainable and productive advancements in animal breeding.

REFERENCES

1. Hayes BJ, Bowman PJ, Chamberlain AJ, Goddard ME. Invited review: genomic selection in dairy cattle: progress and challenges. *J Dairy Sci.* 2009;92:433-43. <https://doi.org/10.3168/jds.2008-1646>
2. Meuwissen THE, Hayes BJ, Goddard ME. Prediction of total genetic value using genome-wide dense marker maps. *Genetics.* 2001;157:1819-29. <https://doi.org/10.1093/genetics/157.4.1819>
3. Wiggans GR, VanRaden PM, Cooper TA. The genomic evaluation system in the United States: past, present, future. *J Dairy Sci.* 2011;94:3202-11. <https://doi.org/10.3168/jds.2010-3866>
4. Anderson CA, Pettersson FH, Clarke GM, Cardon LR, Morris AP, Zondervan KT. Data quality control in genetic case-control association studies. *Nat Protoc.* 2010;5:1564-73. <https://doi.org/10.1038/nprot.2010.116>
5. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al. PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am J Hum Genet.* 2007;81:559-75. <https://doi.org/10.1086/519795>
6. García-Ruiz A, Cole JB, VanRaden PM, Wiggans GR, Ruiz-López FJ, Van Tassell CP. Changes in genetic selection differentials and generation intervals in US Holstein dairy cattle as a result of genomic selection. *Proc Natl Acad Sci USA.* 2016;113:E3995-4004. <https://doi.org/10.1073/pnas.1519061113>
7. Laurie CC, Doheny KF, Mirel DB, Pugh EW, Bierut LJ, Bhangale T, et al. Quality control and quality assurance in genotypic data for genome-wide association studies. *Genet Epidemiol.*

- 2010;34:591-602. <https://doi.org/10.1002/gepi.20516>
8. Marees AT, de Kluiver H, Stringer S, Vorspan F, Curis E, Marie-Claire C, et al. A tutorial on conducting genome-wide association studies: quality control and statistical analysis. *Int J Methods Psychiatr Res*. 2018;27:e1608. <https://doi.org/10.1002/mpr.1608>
 9. Andrews S. FastQC: a quality control tool for high throughput sequence data [Internet]. Babraham Bioinformatics. 2010 [cited 2024 Sep 9]. <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>
 10. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics*. 2014;30:2114-20. <https://doi.org/10.1093/bioinformatics/btu170>
 11. Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet J*. 2011;17:10-2. <https://doi.org/10.14806/ej.17.1.200>
 12. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*. 2009;25:1754-60. <https://doi.org/10.1093/bioinformatics/btp324>
 13. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods*. 2012;9:357-9. <https://doi.org/10.1038/nmeth.1923>
 14. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The sequence alignment/map format and SAMtools. *Bioinformatics*. 2009;25:2078-9. <https://doi.org/10.1093/bioinformatics/btp352>
 15. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernysky A, et al. The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res*. 2010;20:1297-303. <https://doi.org/10.1101/gr.107524.110>
 16. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, et al. The variant call format and VCFtools. *Bioinformatics*. 2011;27:2156-8. <https://doi.org/10.1093/bioinformatics/btr330>
 17. Garrison E, Marth G. Haplotype-based variant detection from short-read sequencing. arXiv:1207.3907 [Preprint]. 2012 [cited 2024 Sep 9]. <https://doi.org/10.48550/arXiv.1207.3907>
 18. Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Wang L, et al. A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly*. 2012;6:80-92. <https://doi.org/10.4161/fly.19695>
 19. Wang K, Li M, Hakonarson H. ANNOVAR: functional annotation of genetic variants from high-throughput sequencing data. *Nucleic Acids Res*. 2010;38:e164. <https://doi.org/10.1093/nar/gkq603>
 20. Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G, et al. Integrative genomics viewer. *Nat Biotechnol*. 2011;29:24-6. <https://doi.org/10.1038/nbt.1754>
 21. Kennedy GC, Matsuzaki H, Dong S, Liu WM, Huang J, Liu G, et al. Large-scale genotyping of complex DNA. *Nat Biotechnol*. 2003;21:1233-7. <https://doi.org/10.1038/nbt869>
 22. Matsuzaki H, Dong S, Loi H, Di X, Liu G, Hubbell E, et al. Genotyping over 100,000 SNPs on a pair of oligonucleotide arrays. *Nat Methods*. 2004;1:109-11. <https://doi.org/10.1038/nmeth718>
 23. Matukumalli LK, Lawley CT, Schnabel RD, Taylor JF, Allan MF, Heaton MP, et al. Development and characterization of a high density SNP genotyping assay for cattle. *PLOS ONE*. 2009;4:e5350. <https://doi.org/10.1371/journal.pone.0005350>
 24. Illumina. BovineSNP50 Genotyping BeadChip. San Diego, CA: Illumina; 2009. Pub. No.: 370-2007-029.
 25. Illumina. GenomeStudio Software [Internet]. 2010 [cited 2024 Sep 9]. <https://www.illumina.com>

- com/products/by-type/informatics-products/microarray-software/genomestudio.html
26. Pryce JE, Haile-Mariam M, Goddard ME, Hayes BJ. Identification of genomic regions associated with inbreeding depression in Holstein and Jersey dairy cattle. *Genet Sel Evol.* 2014;46:71. <https://doi.org/10.1186/s12711-014-0071-7>
 27. Kim S, Lim B, Cho J, Lee S, Dang CG, Jeon JH, et al. Genome-wide identification of candidate genes for milk production traits in Korean Holstein cattle. *Animals.* 2021;11:1392. <https://doi.org/10.3390/ani11051392>
 28. Verma SS, de Andrade M, Tromp G, Kuivaniemi H, Pugh E, Namjou-Khales B, et al. Imputation and quality control steps for combining multiple genome-wide datasets. *Front Genet.* 2014;5:370. <https://doi.org/10.3389/fgene.2014.00370>
 29. Lee J, Kim Y, Cho E, Cho K, Sa S, Kim Y, et al. Genomic analysis using Bayesian methods under different genotyping platforms in Korean Duroc pigs. *Animals.* 2020;10:752. <https://doi.org/10.3390/ani10050752>
 30. Golden Helix. SNP & Variation Suite™(Version 8.x) [Software] [Internet]. 2024 [cited 2024 Sep 9]. <http://www.goldenhelix.com>
 31. Edwards AWF. G. H. Hardy (1908) and Hardy–Weinberg equilibrium. *Genetics.* 2008;179:1143–50. <https://doi.org/10.1534/genetics.104.92940>
 32. Nielsen E, Slatkin M. An introduction to population genetics: theory and application. Sunderland, MA: Sinauer Associates; 2013.
 33. Mayo O. A century of Hardy–Weinberg equilibrium. *Twin Res Hum Genet.* 2008;11:249–56. <https://doi.org/10.1375/twin.11.3.249>
 34. Graffelman J, Weir BS. Testing for Hardy–Weinberg equilibrium at biallelic genetic markers on the X chromosome. *Heredity.* 2016;116:558–68. <https://doi.org/10.1038/hdy.2016.20>
 35. Hill WG. Understanding and using quantitative genetic variation. *Philos Trans R Soc B Biol.* 2010;365:73–85. <https://doi.org/10.1098/rstb.2009.0203>
 36. Alexander DH, Novembre J, Lange K. Fast model-based estimation of ancestry in unrelated individuals. *Genome Res.* 2009;19:1655–64. <https://doi.org/10.1101/gr.094052.109>
 37. Minelli C, Thompson JR, Abrams KR, Thakkinstian A, Attia J. How should we use information about HWE in the meta-analyses of genetic association studies? *Int J Epidemiol.* 2008;37:136–46. <https://doi.org/10.1093/ije/dym234>
 38. McVean G. A genealogical interpretation of principal components analysis. *PLOS Genet.* 2009;5:e1000686. <https://doi.org/10.1371/journal.pgen.1000686>
 39. Browning BL, Browning SR. A unified approach to genotype imputation and haplotype-phase inference for large data sets of trios and unrelated individuals. *Am J Hum Genet.* 2009;84:210–23. <https://doi.org/10.1016/j.ajhg.2009.01.005>
 40. Sargolzaei M, Chesnais JP, Schenkel FS. A new approach for efficient genotype imputation using information from relatives. *BMC Genomics.* 2014;15:478. <https://doi.org/10.1186/1471-2164-15-478>
 41. Pearson TA, Manolio TA. How to interpret a genome-wide association study. *JAMA.* 2008;299:1335–44. <https://doi.org/10.1001/jama.299.11.1335>
 42. Kohavi R. A study of cross-validation and bootstrap for accuracy estimation and model selection. In: *Proceedings of the 14th International Joint Conference on Artificial Intelligence (IJCAI)*; 1995; Montreal. p. 1137–43.
 43. Meuwissen T, Hayes B, Goddard M. Genomic selection: a paradigm shift in animal breeding. *Anim Front.* 2016;6:6–14. <https://doi.org/10.2527/af.2016-0002>
 44. Zhang Z, Zhang Q, Ding X. Advances in genomic selection in domestic animals. *Chin Sci Bull.* 2011;56:2655–63. <https://doi.org/10.1007/s11434-011-4632-7>

45. Berry DP, Kearney JF. Imputation of genotypes from low- to high-density genotyping platforms and implications for genomic selection. *Animal*. 2011;5:1162-9. <https://doi.org/10.1017/S1751731111000309>
46. Howie BN, Donnelly P, Marchini J. A flexible and accurate genotype imputation method for the next generation of genome-wide association studies. *PLOS Genet*. 2009;5:e1000529. <https://doi.org/10.1371/journal.pgen.1000529>
47. Howie B, Fuchsberger C, Stephens M, Marchini J, Abecasis GR. Fast and accurate genotype imputation in genome-wide association studies through pre-phasing. *Nat Genet*. 2012;44:955-9. <https://doi.org/10.1038/ng.2354>
48. VanRaden PM. Symposium review: how to implement genomic selection. *J Dairy Sci*. 2020;103:5291-301. <https://doi.org/10.3168/jds.2019-17684>
49. Habier D, Fernando RL, Kizilkaya K, Garrick DJ. Extension of the Bayesian alphabet for genomic selection. *BMC Bioinformatics*. 2011;12:186. <https://doi.org/10.1186/1471-2105-12-186>
50. Misztal I, Legarra A, Aguilar I. Computing procedures for genetic evaluation including phenotypic, full pedigree, and genomic information. *J Dairy Sci*. 2009;92:4648-55. <https://doi.org/10.3168/jds.2009-2064>